

Protein-Ligand Binding Potential of Mean Force Calculations with Hamiltonian Replica Exchange on Alchemical Interaction Grids

David D. L. Minh*

*Department of Chemistry, Illinois Institute of Technology, Chicago, Illinois 60616, USA and
Department of Chemistry, Duke University, Durham, North Carolina 27708, USA*

(Dated: July 19, 2015)

A binding potential of mean force (BPMF) is a free energy of noncovalent association in which one binding partner is flexible and the other is rigid. Expanding on previous work with host-guest systems, I have developed a method to calculate BPMFs for protein-ligand systems. The method is based on replica exchange sampling from multiple thermodynamic states at different temperatures and protein-ligand interaction strengths. Protein-ligand interactions are represented by interpolating precomputed electrostatic and van der Waals grids. Using a simple estimator for thermodynamic length, thermodynamic states are initialized at approximately equal intervals. The method is demonstrated on the Astex diverse set, a database of 85 protein-ligand complexes relevant to pharmacy or agriculture. Fifteen independent simulations of each complex were started using poses from crystallography, docking, or the lowest-energy pose observed in the other simulations. Benchmark simulations completed within three days on a single processor. Overall, protocols initialized using the thermodynamic length estimator were system-specific, robust, and led to approximately even replica exchange acceptance probabilities between neighboring states. In most systems, the standard deviation of the BPMF converges to within $5 k_B T$. Even with low variance, however, the mean BPMF was sometimes dependent on starting conditions, implying inadequate sampling. Within the thermodynamic cycle, free energies estimated based on multiple intermediate states were more precise, and those estimated by single-step perturbation were less precise. The results demonstrate that the method is promising, but that ligand pose sampling and phase space overlap can sometimes prevent precise BPMF estimation. The software used to perform these calculations, Alchemical Grid Dock (ALGDock), is available under the open-source MIT license at <https://github.com/ccbatit/algdock/>.

INTRODUCTION

Fast and accurate predictions of binding free energies between proteins and small organic ligands would have significant impact on designing drugs [1–3] and other modulators of biological processes. The clear relevance of protein-ligand binding affinity prediction in chemical biology and drug discovery has inspired a vast array of physics-based methods (for a broad review, see [4]), each with a different trade-off between computational accuracy and speed.

On one extreme, molecular docking focuses on speed. Docking algorithms are designed to quickly obtain plausible configurations of a protein-ligand complex. Scoring functions are then used to rank one configuration versus another. Docking programs are commonly assessed by their ability to redock ligands into crystallographic structures from which they have been removed. Comparative studies [5, 6] and blinded exercises [7] consistently show that docking methods are adept at generating the native pose but are less competent at giving it the highest rank. In this context, it is not surprising that docking scores are poorly correlated with binding free energies [5, 6, 8, 9].

In contrast, alchemical pathway methods are based on rigorous statistical mechanics. The methods involve sampling from a series of possibly nonphysical thermodynamic states in between end-states where the receptor and ligand are bound and unbound. In the unbound

state, the receptor-ligand nonbonded interaction terms may be switched off or the species may be physically separated. In accordance to the established statistical mechanics of noncovalent binding [10], the receptor is usually allowed full flexibility. Unfortunately, sampling a fully flexible complex from multiple statistical distributions along an alchemical pathway generally requires substantial computing resources; it has been suggested that most published studies are not fully converged [11]! Even sampling the ligand binding pose is challenging [11]. To bypass this difficulty, most pathway calculations pursue relative binding free energies between similar molecules and assume that the binding mode does not change. In absolute binding free energy calculations, ligand sampling issues are usually alleviated by confining the molecule to a specific pose [12, 13]. Within this restricted range of problems, alchemical pathway calculations are amassing a growing track record of accurate prediction (e.g. [14–22]). *The success of these methods suggests that docking may be substantially improved by incorporating rigorous statistical mechanics.*

Implicit ligand theory (ILT) [23], a recently derived statistical mechanics framework for noncovalent association, has the potential to inspire new methods that combine the speed of docking and the rigor of alchemical pathway methods. ILT formally separates receptor and ligand sampling into two distinct stages. Because the first stage of receptor sampling does not require a lig-

and, receptor conformations can be sampled once and used with many different ligands. In contrast, conventional alchemical pathway methods require thorough receptor sampling for every receptor-ligand pair. The second stage is to calculate the binding potential of mean force (BPMF) - the binding free energy between a ligand and a *rigid* receptor configuration (Eq. 2) - for each receptor configuration. Finally, the standard binding free energy is an exponential average of BPMFs (Eq. 3).

While multiple BPMFs are required to estimate a standard binding free energy, the overall calculation should require less computer time than conventional methods because BPMFs are easier to estimate than binding free energies with a fully flexible receptor. The key reason for this speedup is that a BPMF calculation only requires sampling of the ligand, which usually has many fewer degrees of freedom than the complex. Furthermore, nonbonded interactions between a rigid receptor and ligand can be treated by interpolating precomputed three-dimensional grids, a strategy first developed for docking [24]. Once the grid is stored, calculation time no longer depends on the size of the receptor. In contrast, conventional alchemical pathway methods require frequent force evaluation between flexible receptor atoms. As the number of pairwise interactions scales as $O(N^2)$ with N receptor atoms (neglecting cutoffs), the relative efficiency of ILT-based methods will be more pronounced as receptor size increases.

In the first paper on ILT, BPMFs were estimated for simple host-guest systems [23]. Here, the focus is on precise BPMF estimation for protein-ligand systems. Mobley *et al.* [14], Ucisik *et al.* [25] previously computed binding free energies between simple ligands and a rigid protein, T4 lysozyme. In contrast, the present work involves more diverse systems. Mobley *et al.* [14] were equally rigorous, but did not implement specialized methods to significantly speed their calculations compared to flexible-receptor calculations. Ucisik *et al.* [25] developed a fast method based on a number of approximations. As in other BPMF [23] and standard binding free energy [22, 26, 27] calculations, Hamiltonian replica exchange is applied. The main methodological differences between the current and previous work are the adaptive initialization of thermodynamic states and the use of linearly-scaled interaction grids [24] with transformation-based smoothing [28]; these will be discussed in detail in the section on *Theory and Methods*. The method is demonstrated on the Astex diverse set [29], a curated database of 85 high-quality protein-ligand complexes of pharmaceutical or agrochemical interest.

The algorithm is implemented in a new software package, Alchemical Grid Dock (AlGDock), a python module based on the Molecular Modeling Toolkit (MMTK) 2.7.8 [30]. AlGDock is available under the open-source MIT license at <https://github.com/ccbatit/algdock/>.

THEORY AND METHODS

This section reviews implicit ligand theory, details the algorithms in AlGDock, and describes the setup of the BPMF calculations on the Astex diverse set.

Implicit Ligand Theory

For the noncovalent association between a receptor R and ligand L to form a complex RL , $R + L \rightleftharpoons RL$, the standard binding free energy is,

$$\Delta G^\circ = -\beta^{-1} \ln \left(\frac{C^\circ C_{RL}}{C_R C_L} \right), \quad (1)$$

where $\beta = (k_B T)^{-1}$ is the inverse of Boltzmann's constant times the temperature, C° is the standard concentration (1 M = 1/1660 Å³), and C_X is the equilibrium concentration of species $X \in \{R, L, RL\}$. Activities have been assumed to be unity, a reasonable approximation in the limit of low concentrations.

Coordinates of the complex, r_{RL} , are partitioned into receptor (r_R) and ligand internal (r_L) and external (ξ_L) coordinates. Based on this partitioning, the interaction energy is defined as $\Psi(r_{RL}) = \mathcal{U}(r_{RL}) - \mathcal{U}(r_R) - \mathcal{U}(r_L)$, where $\mathcal{U}(r) = U(r) + W(r)$ is an effective potential energy that includes the gas-phase potential energy $U(r)$ and solvation free energy $W(r)$ [10]. A BPMF is an exponential average of interaction energies over ligand coordinates in the binding site [23],

$$B(r_R) = -\beta^{-1} \ln \left(\frac{\int I_\xi e^{-\beta \Psi(r_{RL})} e^{-\beta \mathcal{U}(r_L)} dr_L d\xi_L}{\int I_\xi e^{-\beta \mathcal{U}(r_L)} dr_L d\xi_L} \right) \quad (2)$$

where $I_\xi \equiv I(\xi_L)$ is an indicator function that takes values between 0 and 1 and specifies whether the receptor and ligand are bound or not.

According to ILT, the standard binding free energy ΔG° is related to an exponential average of BPMFs over Boltzmann-distributed receptor configurations r_R ,

$$\Delta G^\circ = -\beta^{-1} \ln \left(\frac{\int e^{-\beta [B(r_R) + \mathcal{U}(r_R)]} dr_R}{\int e^{-\beta \mathcal{U}(r_R)} dr_R} \right) - \beta^{-1} \ln \left(\frac{\Omega C^\circ}{8\pi^2} \right) \quad (3)$$

where $\Omega = \int I_\xi(\xi_L) d\xi_L$ is the volume of the binding site.

Thermodynamic Cycle

A BPMF can also be expressed as a ratio of partition functions [23],

$$B(r_R) = -\beta^{-1} \ln \left(\frac{\int I_\xi e^{-\beta \mathcal{U}(r_{RL})} dr_L d\xi_L}{\int I_\xi e^{-\beta [\mathcal{U}(r_L) + \mathcal{U}(r_R)]} dr_L d\xi_L} \right) \quad (4)$$

In AlGDock, BPMFs are estimated by completing the thermodynamic cycle shown in Figure 1.

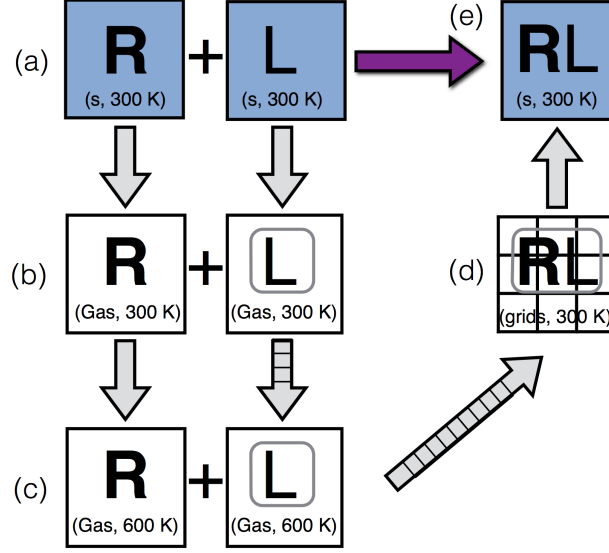


FIG. 1. **Thermodynamic cycle for BPMFs: Schematic.** Milestone thermodynamic states are labeled with letters in parentheses. Counter-clockwise from the top left: the non-interacting receptor and ligand are desolvated (a&b), the temperature is increased to 600 K (b&c), receptor-ligand interaction grids are scaled from 0 to 1 and the temperature is decreased to 300 K (c&d), and the complex is solvated and grid energies replaced with direct calculations (d&e). Arrows with orthogonal lines indicate multiple intermediate thermodynamic states. For BPMF calculations, configurations are sampled from thermodynamic states with the rounded boxes and from their intermediates.

Mathematical expressions for individual free energy differences within the cycle are given in Table I. The sum of all the free energy components in Table I is $\beta B(r_R)$.

TABLE I. Thermodynamic cycle for BPMFs: Expressions		
Milestones	Process	Expression
a&b	Receptor desolvation	$f_{ab,R} = -\ln \frac{e^{-\beta_T U(r_R)}}{e^{-\beta_T U(r_R)}}$
a&b	Ligand desolvation	$f_{ab,L} = -\ln \frac{\int I_\xi e^{-\beta_T U(r_L)} dr_L d\xi_L}{\int I_\xi e^{-\beta_T U(r_L)} dr_L d\xi_L}$
b&c	Receptor heating	$f_{bc,R} = -\ln \frac{e^{-\beta_H U(r_R)}}{e^{-\beta_T U(r_R)}}$
b&c	Ligand heating	$f_{bc,L} = -\ln \frac{\int I_\xi e^{-\beta_H U(r_L)} dr_L d\xi_L}{\int I_\xi e^{-\beta_T U(r_L)} dr_L d\xi_L}$
c&d	Grid scaling/cooling	$f_{cd,o} = -\ln \frac{\int I_\xi e^{-\beta_T U_g(r_{RL})} dr_L d\xi_L}{\int I_\xi e^{-\beta_H [U(r_L) + U(r_R)]} dr_L d\xi_L}$
d&e	Complex solvation	$f_{de} = -\ln \frac{\int I_\xi e^{-\beta_T U(r_{RL})} dr_L d\xi_L}{\int I_\xi e^{-\beta_T U_g(r_{RL})} dr_L d\xi_L}$

Expressions for free energy differences from the thermodynamic cycle in Figure 1. The target temperature $T_T = 300$ K is the basis of $\beta_T = (k_B T_T)^{-1}$. Similarly, the high temperature $T_H = 600$ K is the basis of $\beta_H = (k_B T_H)^{-1}$.

Over the course of the cycle, the receptor-ligand interaction strength is scaled and the temperature is varied. The temperature is varied because high-temperature states enhance transitions between local energetic minima. Because the protocol involves temperature changes, the notation is simplified by using a lower-case u to refer to the *reduced potential energy* [31], a log probability density that incorporates β . For example, at milestone c, the reduced potential energy is $u(r_{RL}) = \beta_H [U(r_L) + U(r_R)]$. Likewise, the reduced free energy difference between two milestones x and y is f_{xy} . Be-

cause converting from reduced to standard potential energies and free energies involves dividing by β , the units of these reduced quantities are specified as $k_B T$.

For thermodynamic states that require I_ξ , the ligand is confined to the binding site using a flat-bottom harmonic potential [23, 26, 27],

$$u_I(d) = \begin{cases} 0 & \text{if } d \leq d_0 \\ \frac{1}{2} \beta k (d - d_0)^2 & \text{if } d > d_0 \end{cases}, \quad (5)$$

where $k = 10000$ kJ/(mol nm²) is the spring constant, d is the distance between the ligand center of mass and the

center of the binding site, and $d_0 = 6.0$ Å is the radius of the binding site. There is no restriction on ligand rotation.

At milestone d, the reduced potential energy is,

$$u_g(r_{RL}) = \beta_T [U(r_L) + U(r_R) + \Psi_g(r_{RL})], \quad (6)$$

where $\Psi_g(r_{RL})$ is the grid interaction energy. The grid interaction energy,

$$\Psi_g(r_{RL}) = \Psi_{PBSA}(r_{RL}) + \Psi_{vdW}(r_{RL}), \quad (7)$$

is based on one electrostatic and two van der Waals grids.

The electrostatic interaction energy $\Psi_{PBSA}(r_{RL})$ is evaluated by multiplying atomic partial charges with the electrostatic potential. The electrostatic potential at each ligand position is obtained by trilinear interpolation of a precomputed grid. The grid is produced by solving the linear Poisson-Boltzmann equation around the minimized receptor molecule using APBS 1.4 [32] with sequential focusing. Coarse grids are at least 1.5 times larger than the range of the receptor molecule in each dimension. Fine grids have the same size as the van der Waals grids, and a spacing of 0.5 Å. Coarse grids use multiple Debye-Huckel boundary conditions, and fine grids use coarse-grid solutions as boundary conditions. Both grids are solved with the following options: a quintic B-spline charge discretization, spline window width of 0.3, protein dielectric of 2.0, solvent dielectric of 80.0, solvent density of 10.0, solvent radius of 1.4 Å, smoothed dielectric and ion-accessibility coefficients, and temperature of 300.0 K.

The van der Waals energy $\Psi_{vdW}(r_{RL})$ is evaluated by an analogous grid-based procedure pioneered by Meng *et al.* [24]. In the AMBER [33] molecular mechanics force field, the receptor-ligand van der Waals interaction energy is represented as, $\sum_{i=1}^{N_{lig}} \sum_{j=1}^{N_{rec}} \left[\frac{A_{ij}}{r_{ij}}^{12} - \frac{B_{ij}}{r_{ij}}^6 \right]$, a double sum over ligand atoms $i \in 1, \dots, N_{lig}$ and receptor atoms $j \in 1, \dots, N_{rec}$. A_{ij} and B_{ij} are the repulsion and attraction parameters and r_{ij} is the distance between atoms i and j . Molecular docking programs usually model these interactions by a geometric mean approximation [24], such that $A_{ij} = \sqrt{A_{ii}}\sqrt{A_{jj}}$ and $B_{ij} = \sqrt{B_{ii}}\sqrt{B_{jj}}$. With this approximation, the receptor-ligand van der Waals interaction energy at a point k is,

$$\Psi_{vdW}(r_{RL}) = \sum_{i=1}^{N_{lig}} \left[\sqrt{A_{ii}}a_k - \sqrt{B_{ii}}b_k \right] \quad (8)$$

$$a_k = \sum_{j=1}^{N_{rec}} \frac{\sqrt{A_{jk}}}{r_{jk}^{12}} \quad (9)$$

$$b_k = \sum_{j=1}^{N_{rec}} \frac{\sqrt{B_{jk}}}{r_{jk}^6}. \quad (10)$$

The values a_k and b_k are precomputed on a three-dimensional grid with a spacing of 0.25 Å. Grids span

at least $10 + 1.5d_{max}$ Å in each dimension surrounding the ligand center of mass in the crystallographic binding mode, where d_{max} is the maximum distance from any ligand atom to its center of mass. For energy and force evaluations, a_k and b_k are approximated at actual ligand atom positions by trilinear interpolation.

Because van der Waals potentials are highly nonlinear, straightforward trilinear interpolation of van der Waals grids provides energies that poorly match directly calculated energies [34]. Hence, van der Waals repulsive grid energies are calculated using a transformation, trilinear interpolation, and inverse transformation procedure [28]: If $f(\mathbf{x}_i)$ is the van der Waals attractive or repulsive energy at \mathbf{x}_i , a position vector on a grid corner, then $g(x_i) = f(x_i)^{-1/m}$ is evaluated for all grid corners. To estimate $f(\mathbf{x})$ at a desired position vector \mathbf{x} , $g(\mathbf{x})$ is obtained by trilinear interpolation of $g(x_i)$ for the eight x_i surrounding x , and $f(x) = g(x)^{-m}$. Previously, when equivalent powers were used for the van der Waals repulsive and attractive grids, $m = 2$ was found to dramatically reduce interpolation error [28]. Here, $m = 2$ is used for the repulsive grid and no transformation for the attractive grid; this combination was found to yield the most accurate van der Waals energies for docked poses.

Alchemical transformations that modulate the strength of interactions between atoms often face an “end-point catastrophe” where free energy changes are numerically unstable [1]. To circumvent this issue, a set of soft Lennard-Jones repulsive and electrostatic grids is introduced between milestones c&d. In these soft grids, the original grid value, v_o , is replaced with $v_{max} \tanh(v_o/v_{max})$. (Gallicchio and Levy [26] also used a hyperbolic tangent energy cap.) For the soft Lennard-Jones repulsive grid, $v_{max} = 10.0$ kJ^{1/2} mol^{-1/2}. To prevent ligands from being “pinned” by electrostatics, the soft electrostatic grid uses as maximum value such that the electrostatic energy is less than or equal to the soft Lennard-Jones repulsive energy for every heavy atom at every grid point. This is established by setting v_{max} for the electrostatic grid to 10 times the minimum ratio of Lennard-Jones and electrostatic scaling factors. The reduced potential energy is switched according to the protocol,

$$\begin{aligned} u_\alpha(r_{RL}) &= \frac{1}{k_B T(\alpha)} [U(r_L) + U(r_R) \\ &\quad + \alpha_{sg}(\alpha) \Psi_{sg}(r_{RL}) + \alpha_g(\alpha) \Psi_g(r_{RL})] \quad (11) \\ \alpha_{sg}(\alpha) &= -(2\alpha - 1)^2 + 1 \\ \alpha_g(\alpha) &= \frac{(2\alpha - 1)^2}{1 + \exp[-1000(\alpha - \frac{1}{2})]} \\ T(\alpha) &= (T_T - T_H)\alpha + T_H \end{aligned}$$

This protocol turns on the soft grids first, and then the unperturbed grids (Figure 2). The potential is consistent with milestone c at $\alpha = 0$ and milestone d at $\alpha = 1$.

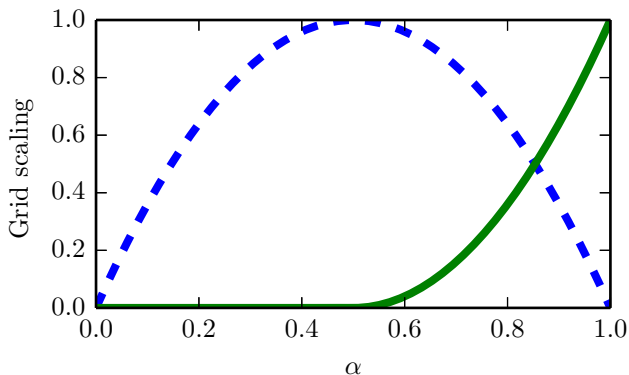


FIG. 2. **Grid scaling between milestones c&d.** α_{sg} (dashed line) and α_g (solid line) as a function of the progress variable α .

Sampling and Estimation

The sampling and estimation strategy in AlGDock is summarized as follows:

1. **Ligand preparation:** the ligand is minimized with 500 conjugate gradient steps. The temperature is ramped from 20 K to 300 K over 30 geometrically spaced simulations of 500 steps each.
2. **Thermodynamic state initialization:** starting from 50 seed configurations, simulations of 1000 steps were used to initialize each thermodynamic state between milestones b&c and c&d.
3. **Hamiltonian replica exchange:** simulations were conducted between milestones b&c and subsequently between c&d.
4. **Postprocessing:** samples from milestones b and d were postprocessed using the generalized Born/surface area (GBSA) water model [35]. For milestone e, grids were replaced by the rigid receptor configuration and pairwise interactions were directly computed.
5. **Estimation:** Free energy differences that sum up to the BPF are estimated.

During thermodynamic state initialization and Hamiltonian replica exchange, the No-U-Turn Sampler (NUTS) [36] is used to propagate the ligand from one configuration to the next. NUTS is a hybrid molecular dynamics (MD)/Markov chain Monte Carlo (MCMC) sampling technique [37] that eliminates the need to select the MD trajectory length between Monte Carlo trials. In other hybrid MD/MCMC techniques, the trajectory length is an adjustable parameter that must be carefully tuned. NUTS automatically truncates MD trajectories when they make a “U-Turn” that brings trajectory end

points closer to one another, or when energy changes are too large. In the present implementation, the change in potential energy or total energy (potential and kinetic energy) is required to be less than 125 kJ/mol. The implementation also truncates trajectories so that the allotted number of MD integration steps is not exceeded.

Another adjustable parameter, the MD time step, is selected using dual averaging - an extension to NUTS [36]. Dual averaging adjusts the time step such that an acceptance statistic H_t becomes close to δ . Parameters for dual averaging were $\delta = 0.6$, $\gamma = 0.05$, $t_0 = 10$, $\kappa = 0.75$, and $\mu = \ln(10\Delta t)$. Because adjusting the time step disturbs the equilibrium distribution, dual averaging is only used during initialization.

To accelerate transitions between binding poses, external coordinate Markov chain Monte Carlo moves are attempted for thermodynamic states between milestones c&d for $\alpha < 0.01$. The external coordinate move consists of:

1. Random rotation. A random quaternion is converted into a matrix that is used to rotate the molecule about its center of mass.
2. Random translation. The magnitude of translation in each dimension is drawn from a Gaussian distribution with a standard deviation of 0.6 Å.

The move is accepted or rejected according to the Metropolis criterion. Moves are not attempted for $\alpha > 0.01$ due to low acceptance probability.

Thermodynamic state initialization

For milestones b&c, the first thermodynamic state ($k = 0$) is at 300 K and the ligand is steadily warmed to 600 K. The first thermodynamic state between milestones c&d depends on whether there is a fully-bound pose available in the binding site. If a pose is available, then the first state is fully bound at 300 K. Otherwise, the first state is fully unbound at 600 K. In the latter situation, 50 randomly selected configurations from milestone c are placed in the binding site at 453 random center-of-mass positions (0.5 positions per Å³) and rotated with 100 different random orientations.

After state k is initialized, state $k + 1$ is initialized as follows:

1. **Parameter selection:** Parameters for state $k + 1$ are selected using a new algorithm designed to separate states at approximately even intervals in thermodynamic length.

Thermodynamic length is a metric of the distance on the manifold of thermodynamic states [38]. For a sequence of states, the statistical error in free energy calculations is minimized and the replica

exchange frequency is nearly maximized when the intermediate states are equidistant in thermodynamic length [39]. Suppose that thermodynamic parameters (e.g. temperature, pressure, grid scaling strengths) are specified by a vector λ with components λ^i . Let $\gamma \equiv \gamma(\alpha)$ describe the dependence of λ on the variable α , such that $\gamma(0)$ is the initial and $\gamma(1)$ is the final thermodynamic state. For microscopic systems, the thermodynamic length is defined by the path integral [39, 40],

$$\mathcal{L} \equiv \int_0^1 \sqrt{\sum_{i,j} \frac{\partial \gamma^i}{\partial \alpha} g(\gamma)_{ij} \frac{\partial \gamma^j}{\partial \alpha}} d\alpha. \quad (12)$$

Given a parameter vector λ , the reduced potential energy is $u_\lambda(x) = \mathcal{U}_\lambda(x)/(k_B T_\lambda)$, where $\mathcal{U}_\lambda(x)$ is the effective potential energy and T_λ is the temperature. The normalized log probability of observing a configuration x is $l_\lambda(x) = -u_\lambda(x) - \ln Z_\lambda$, where $Z_\lambda = \int e^{-u_\lambda(x)} dx$ is the partition function. These quantities are used to define elements of the Fisher information matrix,

$$g(\gamma)_{ij} \equiv \sigma_\lambda^2 [\partial^i l_\lambda, \partial^j l_\lambda], \quad (13)$$

where σ_λ^2 is the covariance in state λ and ∂^i denotes a partial derivative with respect to λ^i . For a protocol in which only one parameter λ^i varies with α , the length is, $\mathcal{L} = \int_0^1 \frac{\partial \gamma^i}{\partial \alpha} \sigma_\lambda [\partial^i l_\lambda] dt$.

Numerical estimates of \mathcal{L} are most accurate when samples are drawn from many states between $0 < \alpha < 1$ [40]. Such exhaustive sampling, however, is unavailable during initialization. A simple approximation for the thermodynamic length when one parameter changes is, $\mathcal{L} = \Delta \lambda^i \sigma_0 [\partial^i l_\lambda]$, where $\Delta \lambda^i$ is the total change in the value of the parameter λ^i and $\sigma_0 [\partial^i l_\lambda]$ is a standard deviation in the initial state. Thus, if one desires \mathcal{L} to be approximately constant between different intermediate stages in a protocol, then the change in parameter should be inversely proportional to $\sigma_0 [\partial^i l_\lambda]$,

$$\Delta \lambda^i = \frac{s}{\sigma_0 [\partial^i l_\lambda]}, \quad (14)$$

where s is an adjustable parameter, the *thermodynamic speed*.

Between milestones b & c , the parameter that varies is the temperature T . As the log probability of a ligand configuration is $l_\lambda = -\frac{U(r_{RL})}{k_B T} - \ln Z_\lambda$, T is incremented by,

$$\Delta \lambda^i = -\frac{s_{bc} k T^2}{\sigma_\lambda [U]}, \quad (15)$$

where $s_{bc} = 0.1$. Between milestones c & d , the log probability of r_{RL} is $l_\lambda = -u_\alpha(r_{RL}) - \ln Z_\lambda$. α is incremented by,

$$\Delta \lambda^i = s_{cd} \left[\left| \frac{d\alpha_{sg}}{d\alpha} \right| \frac{\sigma_\lambda[\Psi_{sg}]}{k_B T(\alpha)} + \left| \frac{d\alpha_g}{d\alpha} \right| \frac{\sigma_\lambda[\Psi_g]}{k_B T(\alpha)} + |T_T - T_H| \frac{\sigma_\lambda[u_\alpha(r_{RL})]}{T(\alpha)} \right]^{-1}, \quad (16)$$

with $s_{cd} = 0.2$.

If the targeted value of the parameter is exceeded (e.g. temperature increased above 600 K), then the targeted value is used.

2. **Seed selection:** 50 configurations from state k are resampled as starting seeds for simulations in state $k + 1$.

Configurations are drawn from state k with weights proportional to $\exp[u_k(x_i) - u_{k+1}(x_i)]$, where $u_k(x) = \mathcal{U}_k(x)/(k_B T_k)$ is the reduced potential energy in state k . In the limit of infinite sampling of state k , resampled configurations would be Boltzmann-distributed in state $k + 1$. With imperfect sampling of state k , resampled configurations approximate the Boltzmann distribution in state $k + 1$.

3. **Sampling:** Simulations of 1000 steps are run from each seed.

NUTS is run with dual averaging to adjust the time step. Each simulation of 1000 steps has its own estimate of the ideal time step. The median time step, restricted to between 0.25 and 2.5 fs, is subsequently used in Hamiltonian replica exchange.

4. **Verification:** The mean replica exchange probability, $\langle p_{acc} \rangle$, is used to verify that thermodynamic states are not too distinct nor similar.

It is estimated by taking the sample mean of p_{acc} (Eq. 17) for every pair of initial samples (at the same time index) from states k and $k + 1$. If $\langle p_{acc} \rangle$ is estimated to be too low (below 0.3), then parameters for state $k + 1$ are reselected with a smaller increment (thermodynamic speed is adjusted by a factor of 4/5) and simulations are repeated. If it is too high (above 0.99), then state k is removed.

Hamiltonian replica exchange

Initialization is followed by Hamiltonian replica exchange [23, 26, 27, 41] between milestones b & c (15 cycles) and c & d (18 cycles). Replica exchange is a Markov chain Monte Carlo move that swaps the configurations of a pair of simulations at different thermodynamic

states. (Equivalently, it may be regarded as swapping the states.) Consider the thermodynamic states a and b with reduced energies u_a and u_b , respectively. If x is the original configuration in state a and y the original configuration in state b , then the acceptance probability,

$$p_{acc} = \min \left[1, e^{-u_a(y) - u_b(x) + u_a(x) + u_b(y)} \right], \quad (17)$$

preserves the Boltzmann distribution in both states.

Typical replica exchange protocols attempt exchanges between pairs of neighboring thermodynamic states, but this restriction is unnecessary. As replica exchange is a type of Gibbs sampling [42], an arbitrary number of attempts can be made between arbitrary pairs of states. In the present simulations, each replica exchange cycle consists of 1000 iterations with 50 NUTS steps, 20 external coordinate MCMC moves (if appropriate), and a sweep of replica exchange. Each sweep of replica exchange includes attempts to swap configurations between pairs of states that are $1, 2, \dots, \min(5, K)$ states apart, where K is the total number of thermodynamic states in the direction.

The statistical inefficiency g is estimated from the integrated autocorrelation time of replica exchange indices, as described [43]. Multiple samples (up to 3 between milestones b&c and 25 between milestones c&d) were saved for every g replica exchange sweeps.

Estimation

BPMFs were estimated according to,

$$B(r_R) = f_{ab,R} + f_{ab,L} + f_{bc,L} + f_{cd} + f_{de}. \quad (18)$$

This expression replaces the sum $f_{bc,R} + f_{cd,o}$ with,

$$f_{cd} = -\ln \frac{\int I_{\xi} e^{-\beta_T [U(r_L) + \Psi_g(r_{RL})]} dr_L d\xi_L}{\int I_{\xi} e^{-\beta_H U(r_L)} dr_L d\xi_L}, \quad (19)$$

which can be evaluated without determining the receptor internal energy $U(r_R)$. The receptor desolvation free energy is estimated by the difference, $f_{ab,R} = \beta_T (U(r_R) - \mathcal{U}(r_R))$.

Other free energy differences are estimated based on equilibrated samples from replica exchange between milestones b&c or c&d. Equilibration is determined by calculating the mean and standard deviation of the energy in milestones b or d, respectively. Systems are considered to be equilibrated once the mean energy is within a standard deviation of the mean energy in the last cycle. $f_{ab,R}$ and f_{de} are estimated by free energy perturbation [44] using configurations drawn from milestones a and d, respectively. $f_{bc,L}$ and f_{cd} are estimated by the multi-state Bennett acceptance ratio [31], which uses potential energies from every replica.

Astex diverse set BPMF calculations

The Astex diverse set is a subset of the protein data bank carefully curated for quality and diversity [29]. The members of the set will be referred to be their protein data bank identifiers, e.g. 1s3v. Each system in the set was set up using AmberTools 14 [45]. Proteins and ions (from protein.mol2) were parameterized using the AMBER ff14SB force field. Zinc parameters were also used for iron, mercury, and magnesium. Non-standard amino acid residues were converted to their standard counterparts. Using the python wrapper to the OpenEye OEChem Toolkit, all crystallographic ligands were parameterized with the Generalized Amber Force Field [46] and AM1BCC charges [47, 48]. The main ligands from ligand.mol were regenerated to lose memory of the crystallographic pose. Noncrystallographic atoms in each receptor were minimized for 2500 steps in gas and 2500 steps in GBSA solvent using NAMD 2.9 [49].

For each system, 15 independent simulations were performed between milestones b&c. Subsequently, simulations were continued between milestones c&d from three different starting points: crystallography, docking, or the lowest-energy pose (according to u_g) sampled in all the other simulations. In most situations, starting from docked poses will be the most practical approach, but the others are a useful comparison.

Docked poses were obtained using the anchor and grow algorithm in UCSF DOCK 6 [50] to redock ligands into their respective proteins. First, sphgen was run with default parameters: dotlim = 0.0, radmax = 4.0 Å, and radmin = 1.4 Å. DOCK was run with 20000 maximum orientations, using internal energy with an exponent of 12, a flexible ligand, an a minimum anchor size of 40. Pruning was performed with clustering, 1000 maximum orientations, a clustering cutoff of 1000, and conformer score cutoff of 25.0. A bump filter was used with max.bumps_anchor = 12 and max.bumps_growth = 12. Final conformations were clustered with a root mean square deviation (RMSD) threshold of 2.0 Å. For two systems, 1t46 and 1n46, no docked poses in the binding site were obtained.

After the starting poses were minimized for 1000 conjugate gradient steps in u_g , the lowest-energy pose was used to initialize milestone d. After the thermodynamic states were initialized, docked poses within the binding site were also used as starting points for each state in replica exchange. The lowest-energy pose was used in milestone d. Higher-energy poses were used in intermediate replicas to fill all available thermodynamic states. If there were more states than docked poses, the lowest-energy pose was duplicated. For simulations starting from a single pose, replicas were not reinitialized at the beginning of replica exchange.

Calculations were performed on the Open Science Grid

[51], Duke shared computing resources, or the Minh group computing cluster at IIT. Each node on the Minh group cluster includes dual AMD Opteron 6376 2.30GHz 16MB cache sixteen-core processors and 64GB DDR3 1600 ECC registered memory. Benchmark calculations were run on the Minh group cluster.

To quantify the precision of the results, we consider the root second moment about a reference value x_r ,

$$\sigma[x, x_r] = \sqrt{\frac{1}{N} \sum_{n=1}^N (x_n - x_r)^2}, \quad (20)$$

for N samples indexed by x_n . When x_r is the mean, then Equation 20 is the standard deviation, which will simply be referred to as $\sigma[x]$. We also consider the root second moment about the minimum, $\sigma[x, x_{min}]$, and about the sample mean at the end of a simulation, $\sigma[x, \bar{x}_{end}]$.

RESULTS

Thermodynamic state initialization

Protocol length

Overall, the present method for thermodynamic state initialization robustly yields protocols that are adapted to specific systems. The large range in the number of thermodynamic states, N_{states} , provides evidence of system-specific adaptation (Figure 3). Between milestones b&c, the average number of states, \bar{N}_{states} ranges from 35.1 to 93.6. Between milestones c&d for simulations initialized from the crystallographic pose, the \bar{N}_{states} ranges from 84.6 to 190.1. For most systems, the standard deviation of the number of states, $\sigma[N_{states}]$ is small relative to the \bar{N}_{states} , indicating that the procedure is robust. The robustness and apparent system-specificity of the protocols implies that a universal thermodynamic protocol may strongly deviate from optimality, e.g. $\langle p_{acc} \rangle$ may be inconsistent across replicas.

In some systems, however, N_{states} is not robust; \bar{N}_{states} is dependent on the starting pose and/or $\sigma[N_{states}]$ is large. Protocols will vary when the initialization stage samples from distinct local minima and estimates different thermodynamic lengths. Between milestones b&c, the $\sigma[N_{states}]$ is larger than 2 in three sets of simulations: 1tz8 (6.1), 1r1h (7.3), and 1vcj (13.6). The larger variance in protocol lengths occurs for different reasons (Figure 4). In 1r1h simulations, the protocols require variable number of states to reach 400 K, after which the slope is similar. For 1tz8 and 1vcj, there are two classes of protocols, indicative of sampling from distinct local minima. With 1tz8, the two protocols are obtained with approximately equal probability. With 1vcj, the short protocol was observed in only 2 of 15 simulations.

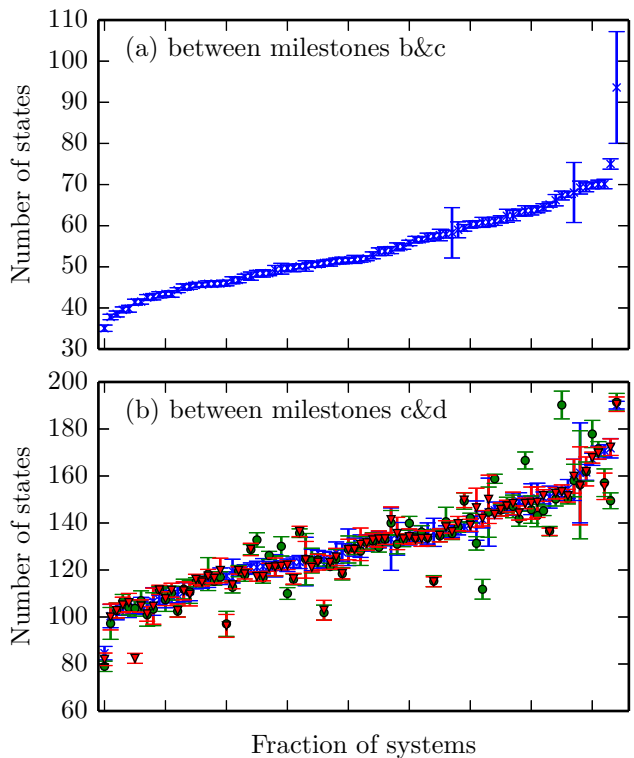


FIG. 3. **Number of thermodynamic states** (a) between milestones b&c, and (b) between milestones c&d. The marker indicates the mean value and error bars the standard deviation of 15 independent simulations. Between milestones c&d, simulations are initialized from the crystallographic (x), lowest-energy docked (circles), or lowest-energy observed (downward-facing triangles) pose. They are ordered by the mean number of states for simulations initialized from the crystallographic pose.

Protocols between milestones c&d are more variable than those between milestones b&c. Interaction grids introduce more possibilities for trapping in local minima. For all three starting points, 117f has the largest variance in N_{states} . The protocols are most similar for the first 40 thermodynamic states and the last 10, but take different paths in between (Figure 5). It is worth noting, however, that high variation in protocols does not necessarily lead to inaccurate results, as will be explained in the following sections.

Time steps

In the vast majority of initialization processes, dual averaging converged to a time step around 1.5 fs (Figure 6 and Table II). This is true both between milestones b&c and c&d, suggesting that the grid does not significantly affect the time step in the vast majority of cases. As a notable exception, the shortest time steps appear

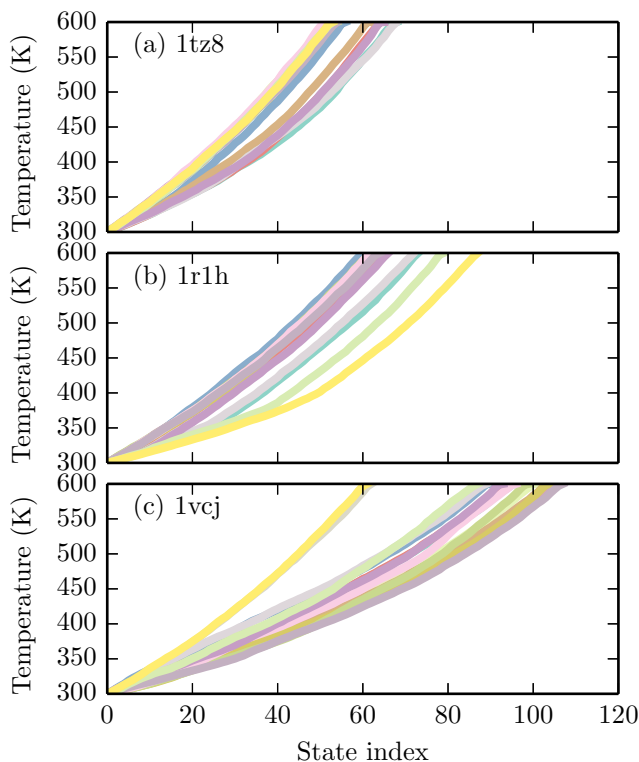


FIG. 4. **Protocols between milestones b&c** for 15 independent simulations of the ligand from (a) 1tz8, (b) 1r1h, and (c) 1vcj.

around $\alpha = 0.5$ between milestones c&d in simulations of 2bm2 and 1n46. In 1n46, there is a steady decrease in time step away from the end points. Because simulations of 1n46 did not start with a docked pose, they are more likely to occupy metastable minima in the grid. Near $\alpha = 0$, the grid has little influence. Near $\alpha = 1$, simulations are dominated by a smaller number of more stable poses. The longest time steps are observed with 1jd0, where most simulations converged to time steps between 1.8 and 2 fs. The ligand in 1jd0 is notable for being small, relatively rigid, and possessing few hydrogen atoms (5). The small number of hydrogen atoms is relevant because unconstrained hydrogen vibrations often limit the maximal time step of molecular dynamics integrators.

TABLE II. **Time step statistics**

Milestones	Mean	Standard Deviation	Min	Max
b&c	1.58	0.11	1.27	2.09
c&d	1.55	0.12	0.87	2.08

Statistics for time steps (in fs) from dual averaging, based on all protocols from all systems.

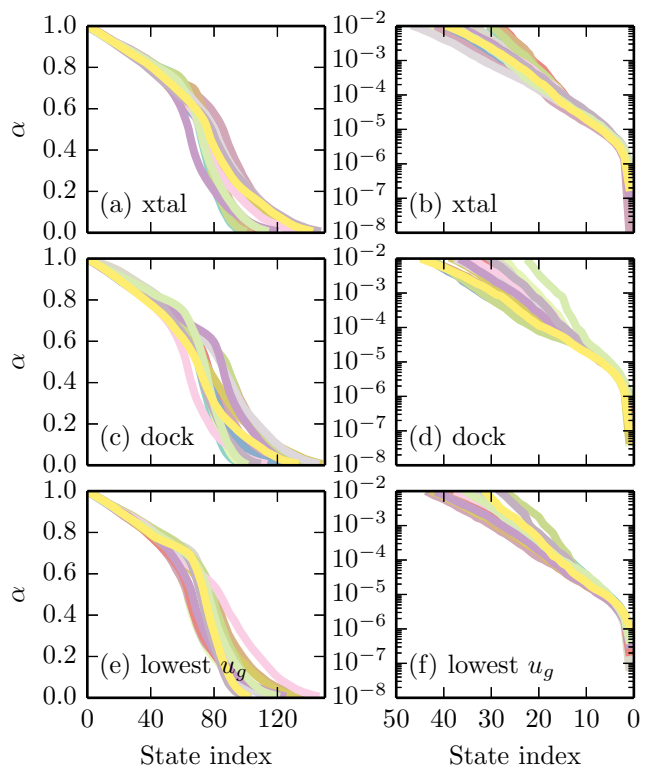


FIG. 5. **Protocols between milestones c&d** for 15 independent simulations of the complex from 1l7f, starting from the crystallographic (a and b), docked (c and d), and lowest-energy observed (e and f) poses. The left panel shows the protocol on a linear scale for $\alpha > 0.01$, indexed such that the first state is at $\alpha = 1.0$. The right panel shows the protocol on a semilog scale for $\alpha < 0.01$, indexed such that the first state is at $\alpha = 0$.

Replica exchange acceptance probabilities

A good replica exchange protocol has a reasonable $\langle p_{acc} \rangle$ between all neighboring states. The key benefit of replica exchange is to spread sampled configurations across a range of different thermodynamic states. A bottleneck in the exchange of configurations across the pair of states can eliminate this benefit of replica exchange; the simulations effectively become two independent sets of simulations. Low exchange probabilities are also indicative of poor phase space overlap, which can limit the convergence of free energy estimates [52].

Because of its importance to efficient simulation, AlGDock includes an estimate of $\langle p_{acc} \rangle$ to verify new thermodynamic states during the initialization step. (The notation \bar{p}_{acc} is used to refer to an *estimator* for $\langle p_{acc} \rangle$). In some simulations, however, the phase space explored during replica exchange can be distinct from that explored during initialization, leading to a substantial change in the observed p_{acc} . Thus, thermodynamic state verification was followed up by estimating $\langle p_{acc} \rangle$ based on equi-

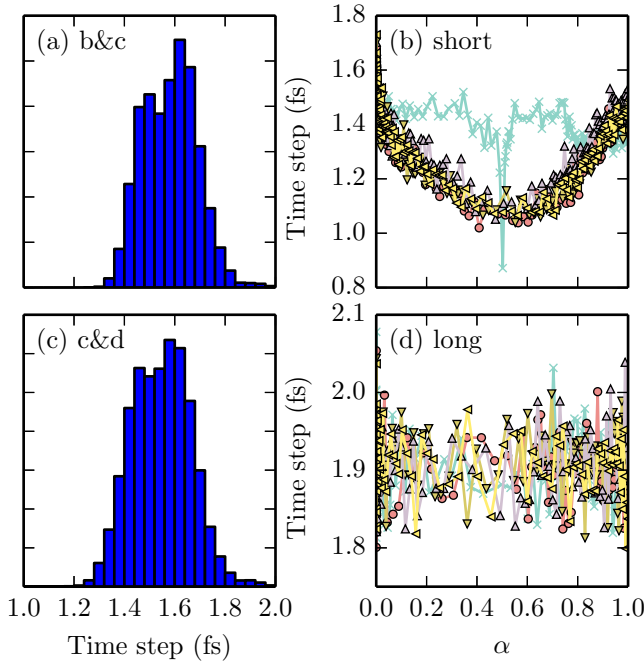


FIG. 6. **Time step statistics.** Statistics for time steps from dual averaging, based on all protocols from all systems. Histograms of time steps between 1.0 and 2.0 fs for (a) between milestones b&c and (c) between milestones c&d. The largest bin count is 10220 between milestones b&c and 63514 between milestones c&d. Five protocols with the (b) shortest and (d) longest observed time steps between milestones c&d. The shortest time steps are observed in 2bm2 (x) and 1n46 (other markers). The longest time steps are observed in 1jd0.

librated samples during replica exchange.

Estimates of $\langle p_{acc} \rangle$ indicate that the initialization protocol does not lead to any simulations with replica exchange bottlenecks (Figure 7 and Table III). Between milestones b&c, \bar{p}_{acc} estimated during replica exchange are high and have low variance. This outcome is consistent with achieving the goal of nearly equal thermodynamic length between adjacent thermodynamic states. It also implies that the configuration spaces explored during initialization and replica exchange are largely the same. Between milestones c&d, the replica exchange rates are also high but there is a larger variance. While most \bar{p}_{acc} are between 0.7 and 1.0, there are a few simulations where \bar{p}_{acc} is much lower; the lowest observed values are during simulations of 2bm2 (0.21) and 1gm8 (0.28). In these simulations, the acceptance probability dips around $\alpha = 0.75$, but are high for most other values of α . These drops indicate that different conformations are explored in replica exchange compared to during thermodynamic state initialization. However, even these low \bar{p}_{acc} are not low enough to be considered a bottleneck; they are expected to allow configurations to pass through the ther-

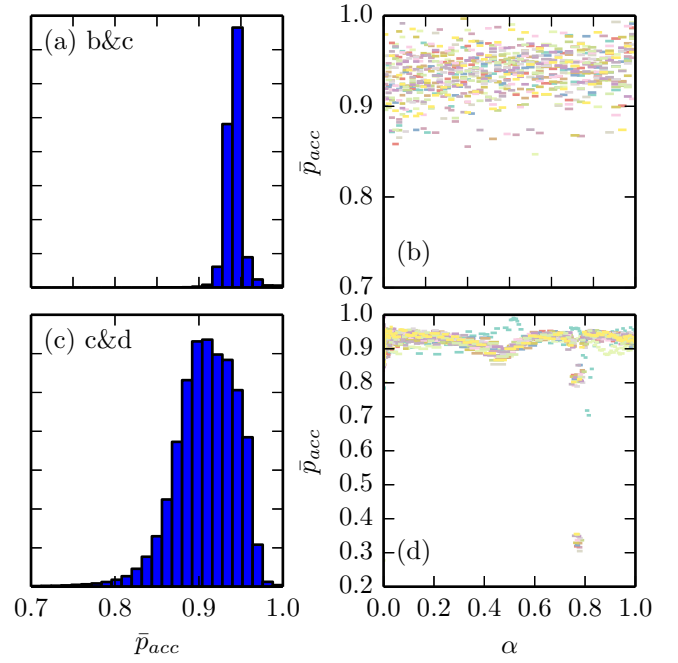


FIG. 7. **Mean acceptance probability statistics.** Statistics for mean acceptance probabilities from replica exchange, based on all simulations. Histograms of \bar{p}_{acc} between 0.7 and 1.0 for (a) between milestones b&c and (c) between milestones c&d. The largest bin count is 35785 between milestones b&c and 63592 between milestones c&d. \bar{p}_{acc} for fifteen protocols with the lowest observed \bar{p}_{acc} (b) between milestones b&c and (d) and between milestones c&d are shown with the a line connecting neighboring states. Between milestones b&c, the temperature increases with α such that $T(\alpha) = (T_H - T_T)\alpha + T_T$. The simulations are from 1s3v (7), 2bm2 (3), 1r55 (3), 1z95 (1), and 1v48 (1). Between milestones c&d, the simulations are from 1v48 (10), 1t40 (3), 2bm2 (1), and 1gm8 (1).

modynamic states several times per cycle.

TABLE III. **Mean acceptance probability statistics**

Milestones	Mean	Standard Deviation	Minimum
b&c	0.94	0.010	0.84
c&d	0.91	0.039	0.21

Statistics for mean acceptance probabilities from replica exchange, \bar{p}_{acc} , based on all simulations.

The observed \bar{p}_{acc} also underscore the point that large variation in protocols is not necessarily problematic. While simulations of 1tz8, 1r1h, and 1vcj have a relatively large variation in N_{states} between milestones b&c, all \bar{p}_{acc} estimated from replica exchange are greater than 0.9 (Figure 8). As is particularly clear for low temperature states of 1r1h, longer protocols have a higher \bar{p}_{acc} . Another bump in \bar{p}_{acc} occurs between the last pair of thermodynamic states, where an even interval in ther-

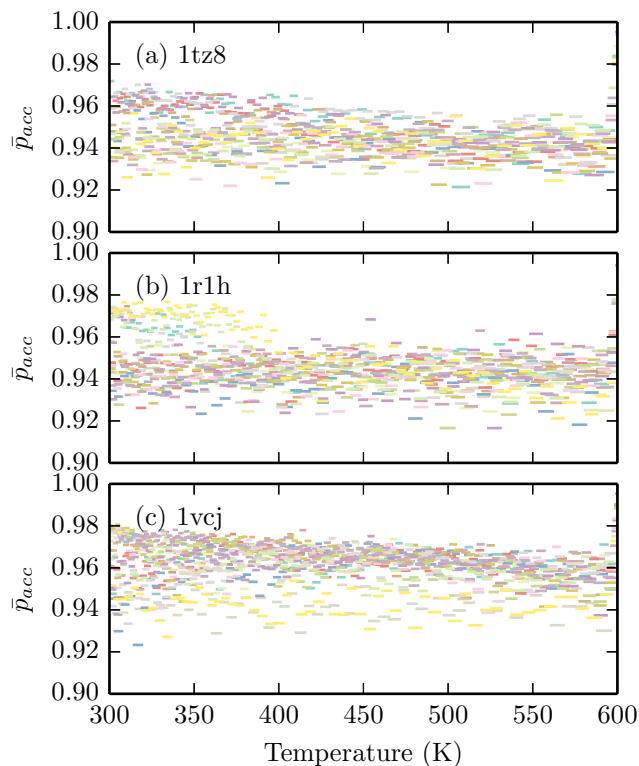


FIG. 8. **Mean acceptance probabilities** for 15 independent simulations of the ligand from (a) 1tz8, (b) 1r1h, and (c) 1vcj between milestones b&c. \bar{p}_{acc} is shown as a line connecting neighboring temperatures.

modynamic length would necessitate going higher than the temperature of interest.

Simulation time

Due to the heterogeneity of computing resources, variability of queue times, and the diverse nature of the systems, calculations took a variable length of time to complete. Some calculations completed in one day, and all were done within a week.

Total CPU time of benchmark simulations spanned a

In many simulations, the phase space sampled at higher α is a subset of that sampled at lower α (e.g. Figure 11). In others, however, the conformational minima for $\alpha \approx 0.5$, where the soft grids are at full strength, are entirely distinct from the minima for $\alpha \approx 1$ (e.g. Figure 10). When there are shifts in important phase space, the thermodynamic states at $\alpha \approx 0.5$ are less beneficial to sampling from milestone d. Nonetheless, phase space overlap between adjacent states is a sufficient condition for precise free energy estimates.

large range, from 13.6 to 71.1 hours (Figure 9). In all simulations, the majority of time was spent in replica exchange between milestones c&d, and the second largest fraction in replica exchange between milestones b&c. Initialization and analysis between milestones c&d consumes a small but noticeable fraction of simulation time, whereas the analogous steps between milestones b&c consume a nearly negligible fraction of simulation time. There is a trend relating the calculation time affiliated with milestones a to c and milestones c to e, but the correlation is imperfect.

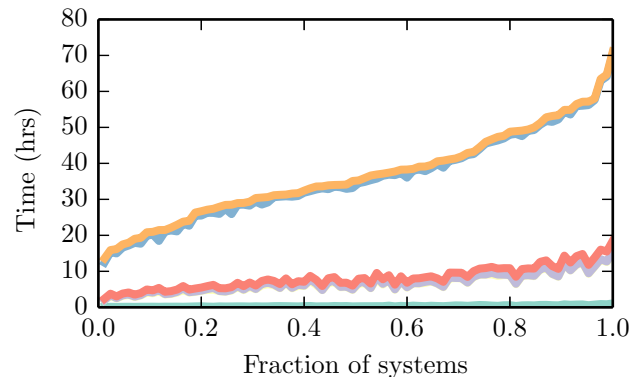


FIG. 9. **Benchmark simulation times** starting from the crystallographic pose. In ascending order, lines depict the cumulative time for initialization, replica exchange, and analysis (running GBSA energy and MBAR calculations), first for milestones a to c and then milestones c to e. Systems are ordered along the x axis by the total simulation time.

Sampling

At milestone c, sampled configurations are uniformly distributed in a sphere. As α increases from 0 to 1 between milestones c&d, the phase space of the ligand is gradually restricted. First, the soft grids prevent ligand atoms from overlapping with receptor atoms. At intermediate values of α , the ligand may assume several poses (e.g. Figure 10). Finally, at milestone d, the ligand usually but does not always sample from a single minimum.

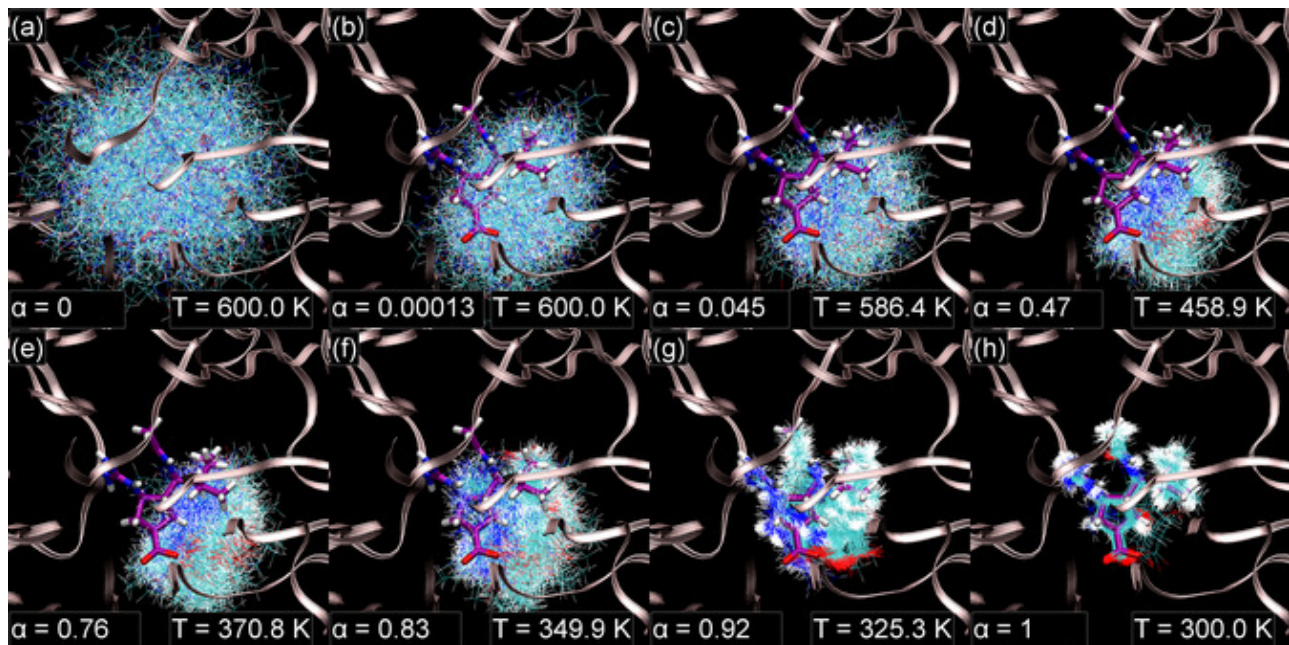


FIG. 10. **Samples from evenly spaced thermodynamic states between milestones c&d**, taken from a representative simulation of 117f starting from docked poses.

At milestone d, most simulations sample from a single minimum (Figure 12), but there are several other situations. These include,

1. Sampling several alternate poses that are (e.g. Figure 12(a)) or are not (e.g. Figure 12(d)) clearly separated;
2. Sampling a clearly defined but sparsely populated alternative pose (e.g. Figure 12(b));
3. Sampling alternate poses that share a common warhead position, but have a floppier tail (e.g. Figure 12(g)).

Native pose identification

A large fraction of simulations include the crystallographic pose among samples from milestone d (Figure 13). Unsurprisingly, it is most often observed in simulations starting from the crystallographic pose. After equilibration, all the simulations still sample a pose with a heavy-atom RMSD less than 0.75 \AA . In contrast, a native pose (defined as having a heavy-atom RMSD less than 2.0 \AA) is observed in 85.6% of simulations starting from the lowest energy pose and 77.7% of simulations starting from docked poses. According to u_g , the lowest-energy observed pose is a native pose in $61/85 = 71.8\%$ of systems. This implies that Boltzmann sampling was able to start at another pose and find the native pose in 13.8% of simulations. Starting docked poses include a native pose in $75/85 = 88.2\%$, and rank the native pose with the lowest u_g (and lowest UCSF DOCK 6 grid score) in $48/85 = 56.5\%$ of systems. This implies that the sampling procedure was able to bring the native pose into milestone d, via Monte Carlo moves or replica exchange, in 21.2% of simulations.

Overall, the GBSA force field is better at identifying the native pose than the grid energy u_g . According to the grid energy, the lowest energy configuration has a RMSD less than 2.0 \AA in 78.5% of simulation starting from the crystallographic pose, 68.9% of simulations starting from the lowest energy pose, and 55.6% of simulations starting from docked poses. According to the GBSA force field, the analogous fractions are 87.1%, 75.1%, and 61.5% of simulations, respectively.

Estimation

Standard deviation of free energy estimates

In most systems, the standard deviation of the BPMF converges to within $5 k_B T$ (Figure 14 and Table IV). Within the thermodynamic cycle, free energies estimated based on multiple intermediate states (f_{bc} and f_{cd}) were more precise, and those estimated by single-step perturbation (f_{ab} and f_{de}) were less precise. Simulating from multiple intermediate states promotes phase space overlap between adjacent states, leading to more precise free

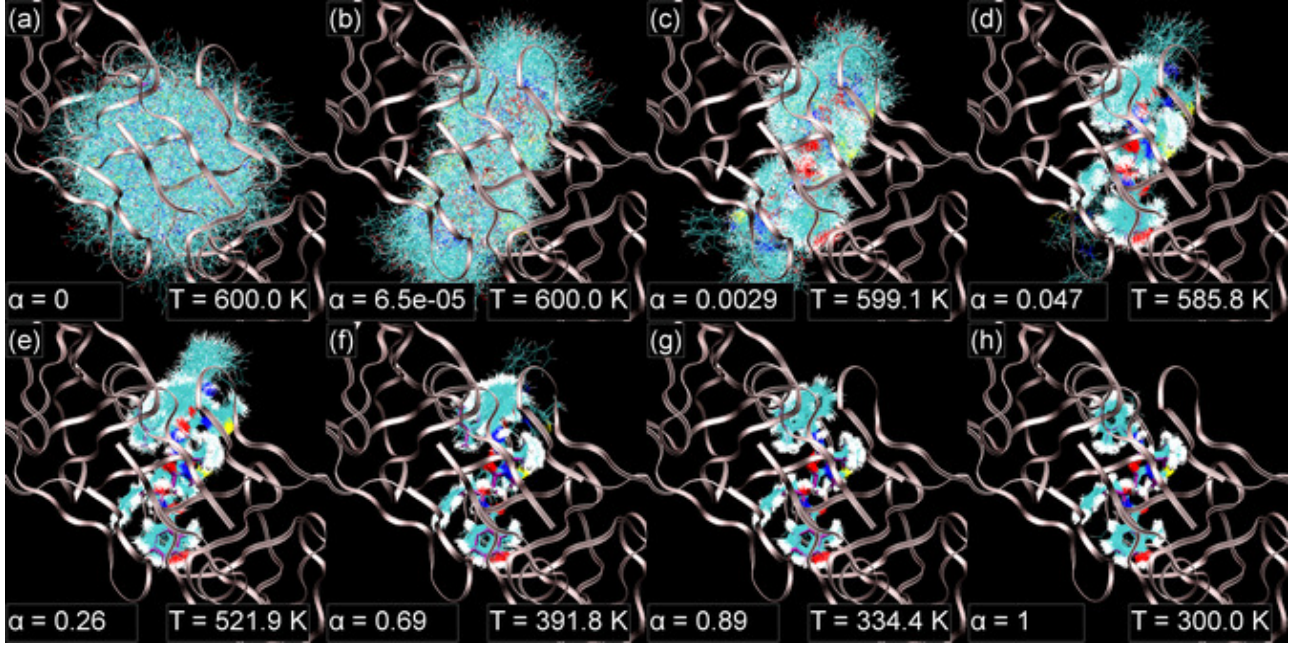


FIG. 11. Samples from evenly spaced thermodynamic states between milestones c&d, taken from a representative simulation of 1kzk starting from docked poses. The protein structure is shown with ribbons and the crystallographic ligand pose is shown with a thick licorice representation and purple carbon atoms. The same illustration scheme is used in Figures 10 and 12. Figures 11, 10, and 12 were generated with VMD [53].

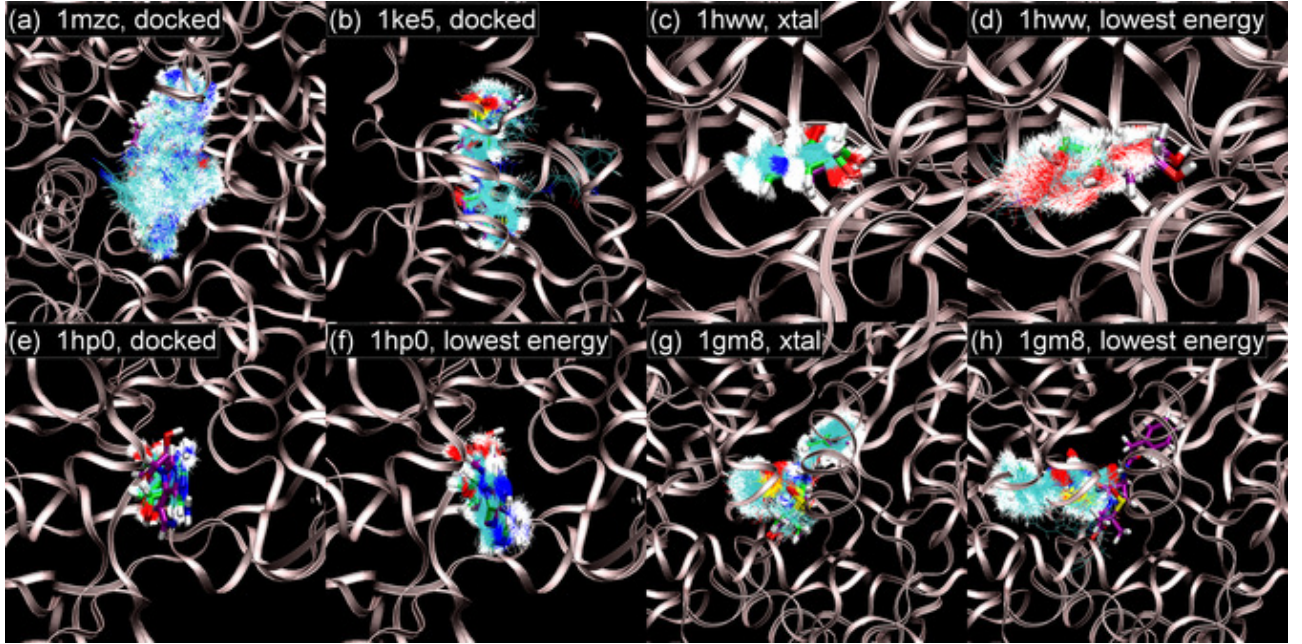


FIG. 12. Samples from milestone d, taken from representative simulations.

energy estimates. Indeed, $\sigma[f_{bc}]$ and $\sigma[f_{cd}]$ are less than $5 k_B T$ except in two systems with no docked poses in the binding site: 1n46 and 1t46. For these simulations, randomly placed ligands end up sampling different local

For systems with poor convergence of BPMPF esti-

minima in independent simulations and result in distinct free energy estimates. Among all free energy differences between adjacent milestones, f_{de} is usually the least precise (Table V).

For systems with poor convergence of BPMPF esti-

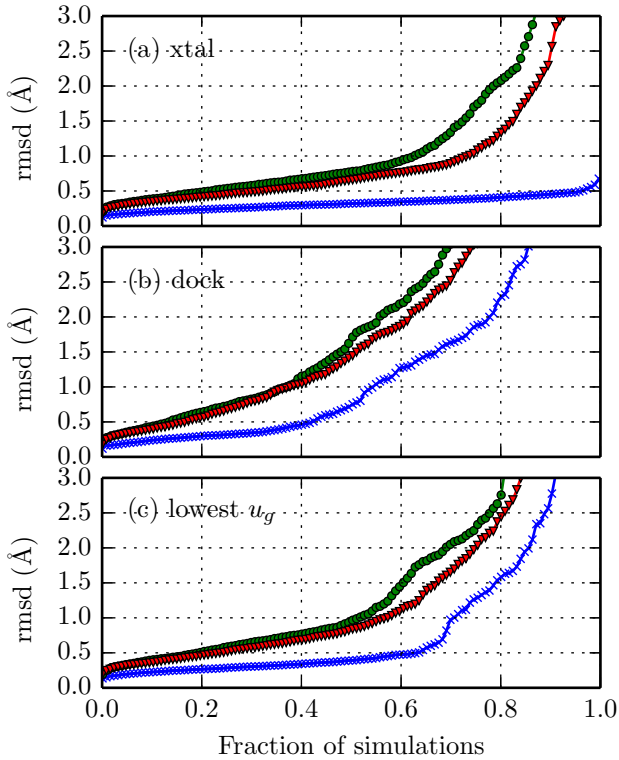


FIG. 13. **Heavy-atom RMSDs from the native state**, for samples from milestone d (after equilibration) from all of the simulations starting from (a) the crystallographic, (b) docked, and (c) lowest-energy observed poses. Symbols denote the lowest observed RMSDs in the simulation (x), or RMSDs of the lowest-energy pose according to u_g (circles) or the GBSA force field (triangles). Each data set is sorted in ascending order. To best emphasize native-like poses, the ordinate is truncated at 3.0 Å.

In some systems, the variance is low but free energy estimates are dependent on the starting poses. When results from all simulations between milestones c&d are grouped together, the standard deviation in many systems no longer converges to within $5 k_B T$. Because converged thermodynamic quantities should not depend on initial conditions, this inconsistency implies that some simulations are not fully converged.

When free energy estimates contradict, it is not completely clear which result is most consistent with the force field. False convergence is quantified based on the minimum value, using $\sigma[x, x_{min}]$. A system is considered to be falsely converged when $\sigma[f]$ is less than $5 k_B T$ but $\sigma[f, f_{min}]$ is greater than $15 k_B T$. Based on this definition, false convergence occurs in a nontrivial fraction of systems (Table IV).

Ultimately, poor convergence and false convergence are caused by incomplete sampling. The most straightforward explanation for incomplete sampling is that a simulation starts and remains trapped in a local min-

imum to improve the precision of free energy estimates. While $\sigma[f, f_{end}]$ usually decreases monotonically with sample size, this trend is not universally true (Figure 15). In some systems there are temporary jumps in $\sigma[f, f_{end}]$ due to instability in determining the equilibrated cycle. More importantly, in other systems, $\sigma[f, f_{end}]$ appears flat, with little or no change with increased sampling.

imum distinct from the global minimum energy. Indeed, between milestones c&d, most falsely converged calculations starting from the crystal structure (4/5) and docked poses (7/11) have no starting pose within 1.0 Å of the lowest-energy observed pose. However, sampling the lowest-energy structure is not a sufficient condition for complete sampling, as other poses may have a lower *free energy*. For simulations starting from the lowest-energy observed pose, four systems appear falsely converged between milestones c&d (formatted by PDB identifier ($\sigma[f]$, $\sigma[f, f_{min}]$) with σ in units of $k_B T$): 1hww (0.353, 16.1), 1r1h (3.95, 15.3), 1hp0 (0.496, 16.7), 1gm8 (0.422, 15.6). For 1hww (Figure 12(c) and (d)) and 1gm8 (Figure 12(g) and (h)), starting from the crystallographic pose leads to a lower estimate of f_{cd} than the lowest-energy pose. For 1hp0 (Figure 12(e) and (f)), a docked pose has the lowest estimated f_{cd} . 1r1h does not appear to sample different poses; the large $\sigma[f, f_{min}]$ is the result of an outlier with a significantly lower f_{cd} estimate than average.

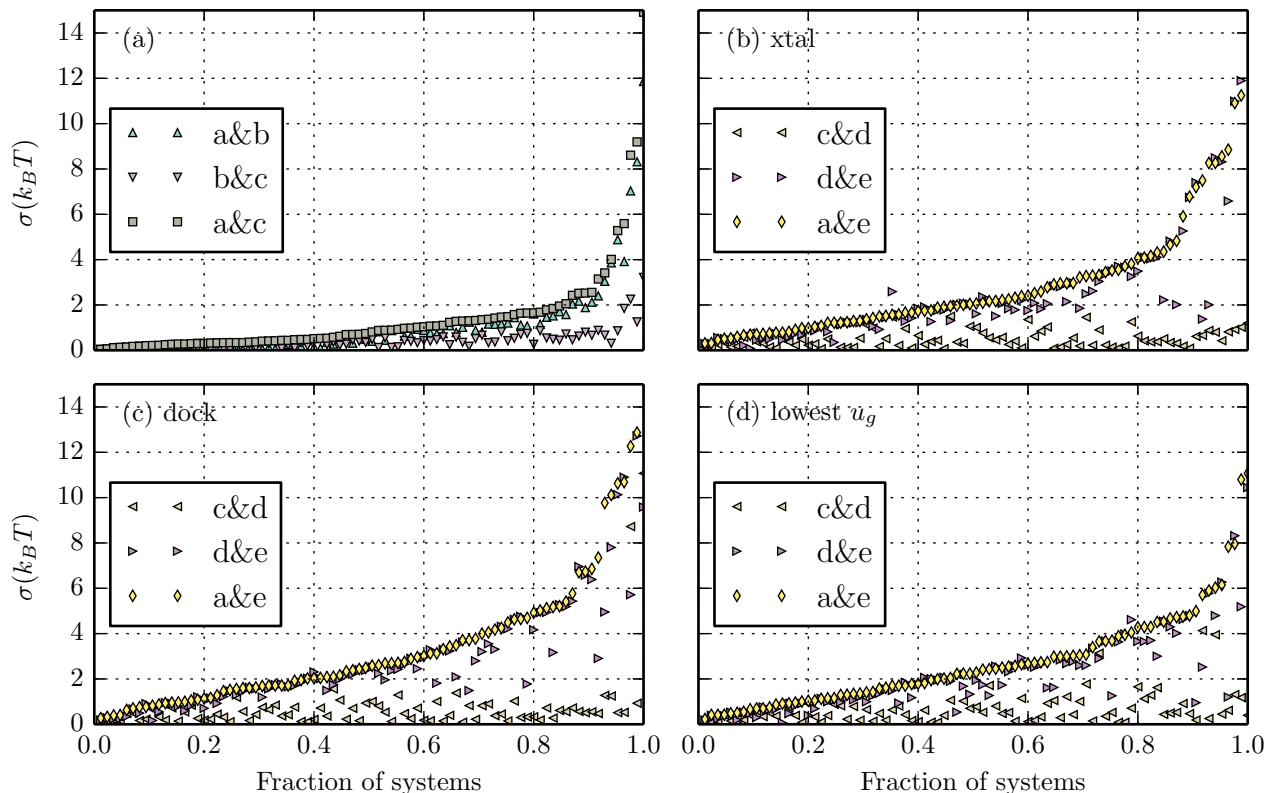


FIG. 14. **Standard deviation of free energy estimates** for (a) 15 independent simulations of the ligand or of the complex starting from (b) the crystallographic, (c) docked, and (d) lowest-energy observed poses. Panel (a) is ordered by the standard deviation of the free energy difference between milestones a&c. The other panels are ordered by the standard deviation of the free energy difference between milestones a&e.

DISCUSSION

The selection of intermediate states between two thermodynamic milestones of interest is a ubiquitous problem in molecular simulation. The usual approach is an iterative trial-and-error process starting from a naive protocol, checking for issues such as replica exchange bottlenecks, and inserting and removing states as necessary. I have developed a simple and robust approach to initialize a series of thermodynamic states based on only a single adjustable parameter, the thermodynamic speed. Due to this automated procedure, I was able to run a large number of simulations on a diverse array of protein-ligand complexes. In the vast majority of cases, the protocols yielded consistent replica exchange rates across neighboring thermodynamic states without further fine-tuning. The described approach to trailblazing thermodynamic state space may find use in other classes of simulations.

Replica exchange calculations in this present study include more thermodynamic states than most published molecular simulations. Conventional wisdom about replica exchange is that an optimal number of replicas will maximize efficiency. With too few replicas, there is

limited phase space overlap between neighbors and exchange rates are vanishingly small. With too many replicas, metrics of replica exchange efficiency, such as a mean round-trip time, diminish. However, in a recent study involving extensive simulation of several distinct processes, it was found that if there are no bottlenecks, the number of states has little impact on the convergence of free energy estimates (the manuscript is under preparation). Hence, I chose to include a large number of states to allow for the possibility that later sampling will explore different regions of phase space and reduce the exchange rate between neighboring states.

While useful, the described thermodynamic state initialization process remains imperfect. Replica exchange is particularly beneficial when the important phase space of a thermodynamic states is a subset of the important phase space of another. The present procedure is limited to the variation of a single thermodynamic parameter between milestones, and can do little to promote the subset relationship. Future improvements could accommodate varying multiple parameters (e.g. separate parameters for the temperature, van der Waals grids, and electrostatic grids) in between thermodynamic states of interest, while

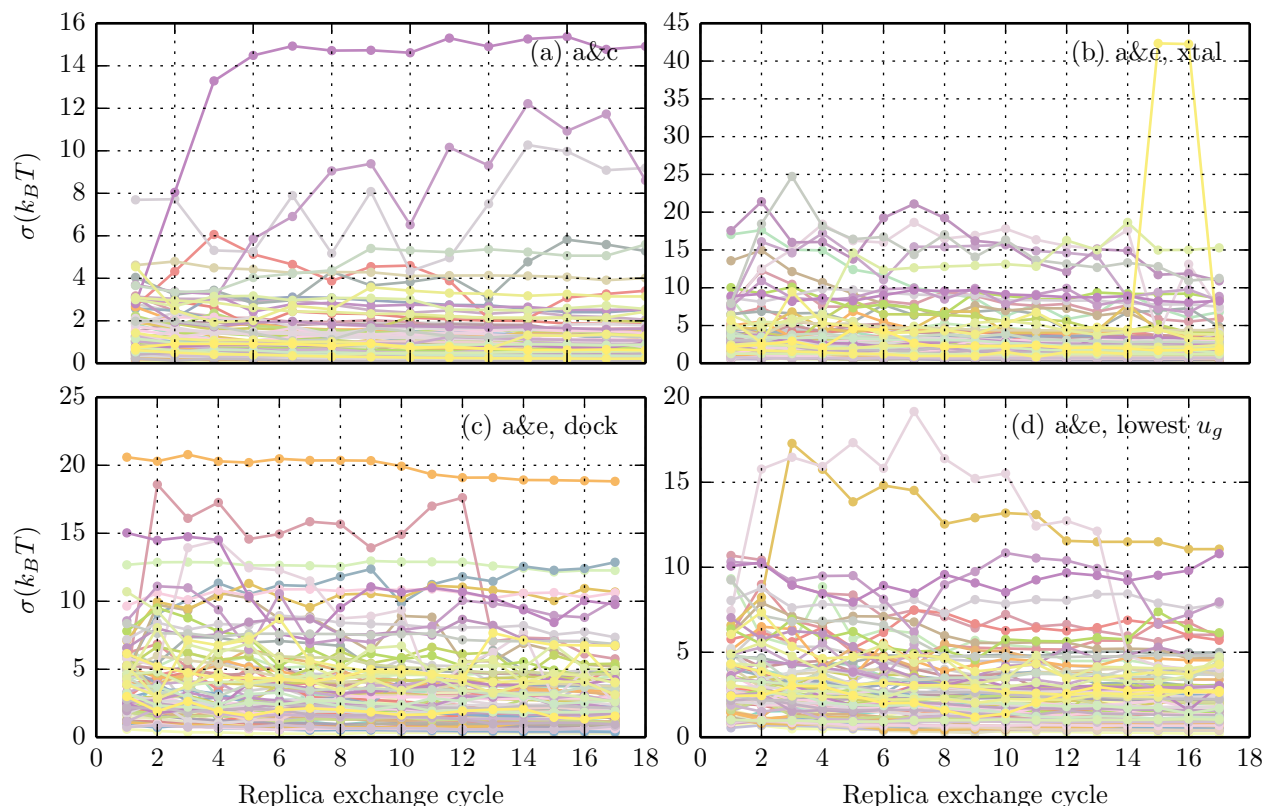


FIG. 15. **Standard deviation of free energy estimates as a function of cycle:** $\sigma[f_{ac}, \bar{f}_{ac,end}]$ based on 15 independent simulations of (a) the ligand or $\sigma[f_{ae}, \bar{f}_{ae,end}]$ based on the ligand simulations and 15 independent simulations of the complex starting from the (b) crystallographic, (c) docked, and (d) lowest-energy observed poses.

maintaining a phase space subset relationship.

While a direct comparison is difficult, the described procedure appears less successful than the best docking programs at identifying the native pose. Hartshorn *et al.* [29] found that with a larger site (10 Å), the top-ranked GOLD solution is within 2 Å of the experimental binding mode in 80.4% of calculations. In contrast, only 56.5% of the present UCSF DOCK 6 calculations were successful, suggesting that setup and parameterization, and perhaps the docking algorithm, could be fine-tuned. Native pose identification based on AIGDock performs similarly to UCSF DOCK 6, suggesting that replica exchange sampling from the Boltzmann distribution does not significantly increase the diversity of poses. As evidenced by the variable performance as a function of starting pose(s), improving Boltzmann sampling schemes should improve native pose identification.

Nonetheless, the sampling and estimation approach is able to attain precise BPMF estimates in the majority of systems. As for the others, it is unsurprising that the same approach does not work equally well in all cases. With a diverse set, ligands and receptors will exhibit distinct levels of complexity and may be driven to bind by different factors. Indeed, the inadequate sampling that

led to convergence issues for several systems is probably not unique to BPMF calculations [11]. Because of the relative simplicity of the calculations, I have been able to conduct more replicates from more distinct starting points than other alchemical standard binding free energy calculations. The sheer number of simulations has allowed us to expose problems that may not be clear from fewer calculations. Our results suggest that Hamiltonian replica exchange with molecular dynamics is not always an adequate approach for sampling distinct ligand binding poses. Specialized sampling methods will likely be required in a more broadly effective protein-ligand BPMF estimation procedure. It is also clear that for some systems, there may be limited phase space overlap between different implicit solvent models. More precise BPMF estimation can be attained by introducing intermediate states or implementing more similar force fields for milestones a&b and d&e.

CONCLUSIONS

I have developed a method to initialize thermodynamic states and estimate BPMFs for protein-ligand systems.

TABLE IV. Number of converged systems

Milestones	Starting point	$< 1.5 k_B T$	$< 5.0 k_B T$	False
a&b	-	69	82	-
b&c	-	82	85	-
a&c	-	64	80	-
c&d	xtal	84	85	5
c&d	dock	82	83	9
c&d	lowest u_g	77	85	4
c&d	all	48	71	-
d&e	xtal	34	76	4
d&e	dock	26	73	8
d&e	lowest u_g	37	80	11
d&e	all	17	63	-
a&e	xtal	28	74	6
a&e	dock	20	68	11
a&e	lowest u_g	27	77	12
a&e	all	15	52	-

Number of systems (out of 85) in which the free energy difference has a standard deviation of less than 1.5 or less than $5.0 k_B T$. False convergence means that the system $\sigma[x] < 5.0 k_B T$, but $\sigma[x, x_{min}] > 15.0 k_B T$. For the xtal, dock, and lowest u_g starting points, the standard deviation is based on 15 independent simulations. The starting point of ‘all’ combines data from the other starting points, and is thus based on 45 independent simulations.

TABLE V. Least precise free energy components

Starting point	Subset	a&b	b&c	c&d	d&e
xtal	all	11	0	5	69
dock	all	10	1	5	69
lowest u_g	all	12	1	9	63
all	all	6	0	20	59
xtal	imprecise	5	0	0	10
dock	imprecise	3	0	2	18
lowest u_g	imprecise	4	0	1	13
all	imprecise	3	0	11	25

The least precise free energy estimates between adjacent milestones. Counts are based on all simulations or the imprecise subset, when the standard deviation of the total BPFM is greater than 5.0.

The largest sources of imprecision are found to be ligand pose sampling and phase space overlap between implicit solvent models. Intriguingly, I found that starting from the lowest-energy pose is not a sufficient condition for converged BPFM estimates.

Future studies may build on the current work to compare BPFMs with molecular docking scores as classifiers of binding. Another useful direction will be to assess the convergence of standard binding free energy estimates

upon increased receptor conformational sampling.

ACKNOWLEDGMENTS

I thank Michael Shirts (University of Virginia), John Chodera (MSKCC), and Trung Hai Nguyen (IIT) for helpful discussions. Peter Eastman (Simbios) and John Chodera assisted with implementing an early version of the code in OpenMM. Michael Sherman (Simbios) made an invaluable suggestion of using precomputed grids. Rob Gardner, Lincoln Bryant, and Balamurugan Desinghu (Open Science Grid) and Tom Milledge (DSCR) provided assistance with computing resources. David Beratan (Duke) was a supportive postdoctoral advisor. OpenEye Scientific Software, Inc. and UCSF provided academic licenses to their software. At the beginning of this project, I was a postdoctoral scholar at Duke, supported by NIH 2P50 GM-067082-06-10. I also spent a month as an OpenMM visiting scholar at Simbios.

REFERENCES

-
- * Electronic Address: david.minh@iit.edu
- [1] J. Michel and J. W. Essex, Journal of Computer-Aided Molecular Design **24**, 639 (2010).
 - [2] J. D. Chodera, D. L. Mobley, M. R. Shirts, R. W. Dixon, K. Branson, and V. S. Pande, Current Opinion in Structural Biology **21**, 150 (2011).
 - [3] D. L. Mobley and P. V. Klimovich, Journal of Chemical Physics **137**, 230901 (2012).
 - [4] M. K. Gilson and H.-X. Zhou, Annual Review of Biophysics and Biomolecular Structure **36**, 21 (2007).
 - [5] G. L. Warren, C. W. Andrews, A.-M. Capelli, B. Clarke, J. LaLonde, M. H. Lambert, M. Lindvall, N. Nevins, S. F. Semus, S. Senger, G. Tedesco, I. D. Wall, J. M. Woolven, C. E. Peishoff, and M. S. Head, Journal of Medicinal Chemistry **49**, 5912 (2006).
 - [6] R. Kim and J. Skolnick, Journal of Computational Chemistry **29**, 1316 (2008).
 - [7] K. L. Damm-Ganamet, R. D. Smith, J. B. Dunbar, J. A. Stuckey, and H. A. Carlson, Journal of Chemical Information and Modeling **53**, 1853 (2013).
 - [8] R. Wang, Y. Lu, X. Fang, and S. Wang, Journal of Chemical Information and Computer Sciences **44**, 2114 (2004).
 - [9] R. D. Smith, J. B. Dunbar, P. M.-U. Ung, E. X. Esposito, C.-Y. Yang, S. Wang, and H. A. Carlson, Journal of Chemical Information and Modeling **51**, 2115 (2011).
 - [10] M. K. Gilson, J. A. Given, B. L. Bush, and J. A. McCammon, Biophysical Journal **72**, 1047 (1997).
 - [11] D. L. Mobley, Journal of Computer-Aided Molecular Design **26**, 93 (2012).
 - [12] S. Boresch, F. Tettinger, M. Leitgeb, and M. Karplus, Journal of Physical Chemistry B **107**, 9535 (2003).

- [13] D. L. Mobley, J. D. Chodera, and K. A. Dill, *Journal of Chemical Physics* **125**, 84902 (2006).
- [14] D. L. Mobley, A. P. Graves, J. D. Chodera, A. C. McReynolds, B. K. Shoichet, and K. A. Dill, *Journal of Molecular Biology* **371**, 1118 (2007).
- [15] J. Michel and J. W. Essex, *Journal of Medicinal Chemistry* **51**, 6654 (2008).
- [16] S. E. Boyce, D. L. Mobley, G. J. Rocklin, A. P. Graves, K. A. Dill, and B. K. Shoichet, *Journal of Molecular Biology* **394**, 747 (2009).
- [17] Y. Deng and B. Roux, *Journal of Physical Chemistry B* **113**, 2234 (2009).
- [18] X. Ge and B. Roux, *Journal of Physical Chemistry B* **114**, 9525 (2010).
- [19] R. D. Malmstrom and S. J. Watowich, *Journal of Chemical Information and Modeling* **51**, 1648 (2011).
- [20] L. Wang, B. J. Berne, and R. A. Friesner, *Proceedings of the National Academy of Sciences of the United States of America* **109**, 1937 (2012).
- [21] S. Zhu, S. M. Travis, and A. H. Elcock, *Journal of Chemical Theory and Computation* **9**, 3151 (2013).
- [22] L. Wang, Y. Wu, Y. Deng, B. Kim, L. Pierce, G. Krilov, D. Lupyan, S. Robinson, M. K. Dahlgren, J. Greenwood, D. L. Romero, C. Masse, J. L. Knight, T. Steinbrecher, T. Beuming, W. Damm, E. Harder, W. Sherman, M. Brewer, R. Wester, M. Murcko, L. Frye, R. Farid, T. Lin, D. L. Mobley, W. L. Jorgensen, B. J. Berne, R. a. Friesner, and R. Abel, *Journal of the American Chemical Society* **137**, 2695 (2015).
- [23] D. D. L. Minh, *Journal of Chemical Physics* **137**, 104106 (2012).
- [24] E. C. Meng, B. K. Shoichet, and I. D. Kuntz, *Journal of Computational Chemistry* **13**, 505 (1992).
- [25] M. N. Ucisik, Z. Zheng, J. C. Faver, and K. M. Merz, *Journal of Chemical Theory and Computation* **10**, 1314 (2014).
- [26] E. Gallicchio and R. M. Levy, *Journal of Computer-Aided Molecular Design* **26**, 505 (2012).
- [27] K. Wang, J. D. Chodera, Y. Yang, and M. R. Shirts, *Journal of Computer-Aided Molecular Design* **27**, 989 (2013).
- [28] C. M. Venkatachalam, X. Jiang, T. Oldfield, and M. Waldman, *Journal of Molecular Graphics and Modelling* **21**, 289 (2003).
- [29] M. J. Hartshorn, M. L. Verdonk, G. Chessari, S. C. Brewerton, W. T. M. Mooij, P. N. Mortenson, and C. W. Murray, *Journal of Medicinal Chemistry* **50**, 726 (2007).
- [30] K. Hinsin, *Journal of Computational Chemistry* **21**, 79 (2000).
- [31] M. R. Shirts and J. D. Chodera, *Journal of Chemical Physics* **129**, 124105 (2008).
- [32] N. A. Baker, D. Sept, S. Joseph, M. J. Holst, and J. A. McCammon, *Proceedings of the National Academy of Sciences of the United States of America* **98**, 10037 (2001).
- [33] J. W. Ponder and D. A. Case, *Advances in Protein Chemistry* **66**, 27 (2003).
- [34] D. Oberlin and H. A. Scheraga, *Journal of Computational Chemistry* **19**, 71 (1998).
- [35] I. Massova and P. A. Kollman, *Perspectives In Drug Discovery And Design* **18**, 113 (2000).
- [36] M. D. Hoffman and A. Gelman, (2011), arXiv:1111.4246v1.
- [37] S. Duane, A. D. Kennedy, B. J. Pendleton, and D. Roweth, *Physics Letters B* **195**, 216 (1987).
- [38] F. Weinhold, *Journal of Chemical Physics* **63**, 2479 (1975).
- [39] D. K. Shenfeld, H. Xu, M. P. Eastwood, R. O. Dror, and D. E. Shaw, *Physical Review E* **80**, 46705 (2009).
- [40] G. E. Crooks, *Physical Review Letters* **99**, 100602 (2007).
- [41] W. Jiang and B. Roux, *Journal of Chemical Theory and Computation* **6**, 2559 (2010).
- [42] J. D. Chodera and M. R. Shirts, *Journal of Chemical Physics* **135**, 194110 (2011).
- [43] J. D. Chodera, W. C. Swope, J. W. Pitera, C. Seok, and K. A. Dill, *Journal of Chemical Theory and Computation* **3**, 26 (2007).
- [44] R. Zwanzig, *Journal of Chemical Physics* **22**, 1420 (1954).
- [45] D. A. Case, T. A. Darden, T. E. Cheatham III, C. L. Simmerling, J. Wang, R. E. Duke, R. Luo, R. C. Waker, W. Zhang, K. Merz, B. Roberts, S. Hayik, A. E. Roitberg, G. Seabra, J. Swails, A. W. Goetz, I. Kolossváry, K. Wong, F. Paesani, J. Vanicek, R. Wolf, J. Liu, X. Wu, S. Brozell, T. Steinbrecher, H. Gohlke, Q. Cai, X. Ye, J. Wang, M.-J. Hsieh, G. Cui, D. R. Roe, D. Mathews, M. Seetin, R. Salomon-Ferrer, C. Sagui, V. Babin, T. Luchko, S. Gusarov, A. Kovalenko, and P. Kollman, "AMBER," (2012).
- [46] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, and D. A. Case, *Journal of Computational Chemistry* **25**, 1157 (2004).
- [47] A. Jakalian, B. L. Bush, D. B. Jack, and C. I. Bayly, *Journal of Computational Chemistry* **21**, 132 (1999).
- [48] A. Jakalian, D. B. Jack, and C. I. Bayly, *Journal of Computational Chemistry* **23**, 1623 (2002).
- [49] J. C. Phillips, R. Braun, W. Wang, J. Gumbart, E. Tajkhorshid, E. Villa, C. Chipot, R. D. Skeel, L. Kalé, and K. Schulten, *Journal of Computational Chemistry* **26**, 1781 (2005).
- [50] P. Lang, S. Brozell, S. Mukherjee, E. Pettersen, E. C. Meng, V. Thomas, R. C. Rizzo, D. A. Case, T. James, and I. D. Kuntz, *RNA* **15**, 1219 (2009).
- [51] R. Pordes, D. Petravick, B. Kramer, D. Olson, M. Livny, A. Roy, P. Avery, K. Blackburn, T. Wenaus, F. Würthwein, I. Foster, R. Gardner, M. Wilde, A. Blatecky, J. McGee, and R. Quick, *Journal of Physics: Conference Series* **78**, 012057 (2007).
- [52] N. Lu and D. A. Kofke, *Journal of Chemical Physics* **114**, 7303 (2001).
- [53] W. Humphrey, A. Dalke, and K. Schulten, *Journal of Molecular Graphics* **14**, 33 (1996).